

Knowledge Containers tagged to Taxonomies

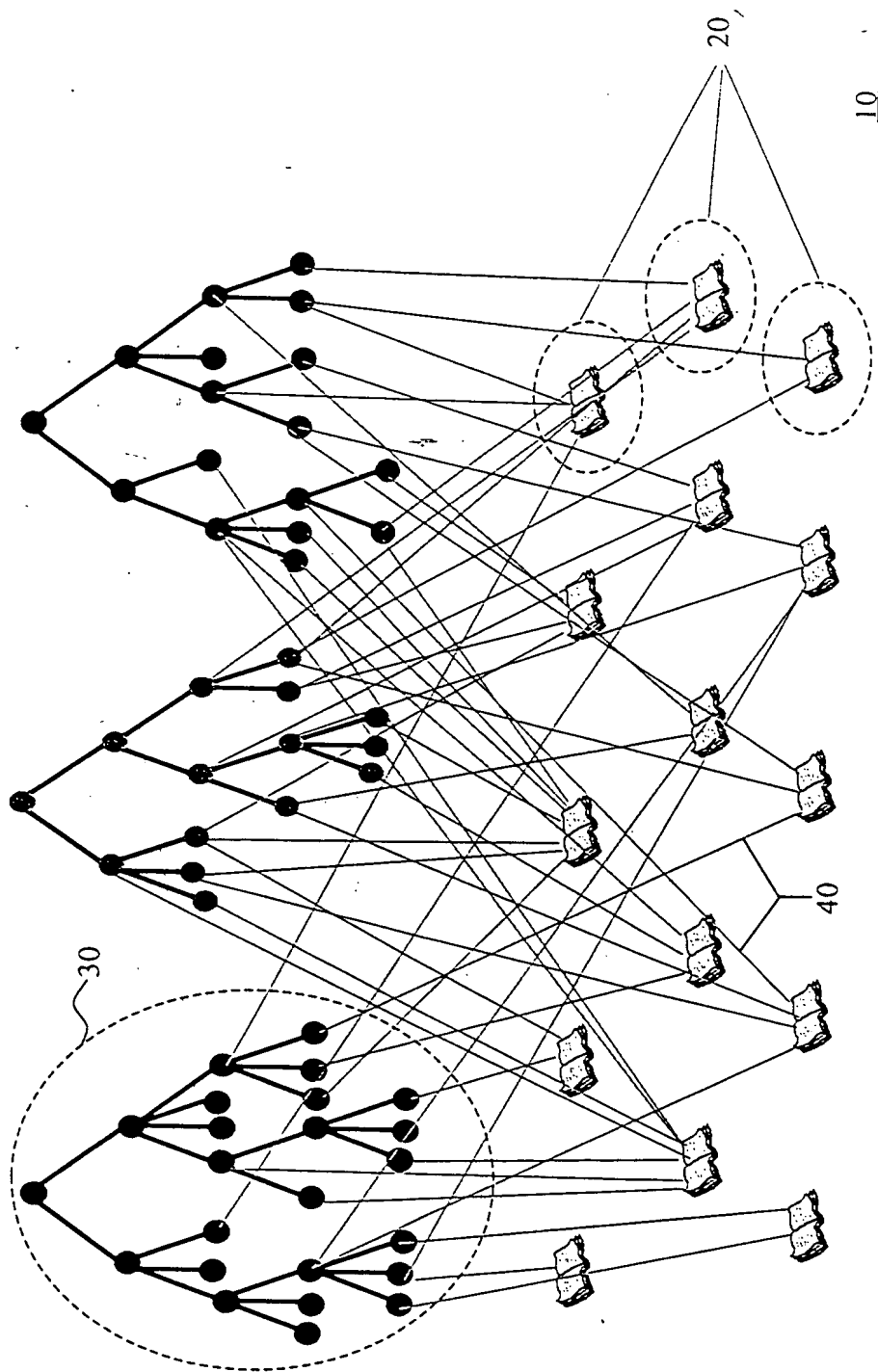


FIG. 1

A Knowledge Container


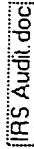

<i>Administrative meta-data</i>	<p>50</p> <p><author>Rev. Bill C. Wurtz</author> <creation date>6/7/89</creation date> <expiration date>12/31/99</expiration date></p>
<i>Taxonomy Tags</i>	<p>60</p> <p>Tax_Audit: 0.92 Tax_Evasion: 0.65 Fraud: 0.45</p>
<i>Marked Content</i>	<p>70</p> <p><P> In 1988, the <org>IRS</org> investigated <person>Scott Huffman</person>, looking for irregularities in income reporting. <P> ... <P>The <term>preliminary charges</term> included ...</p>
<i>Original Content</i>	<p>80</p> <p> </p>
<i>Links</i>	<p>90</p> <p> mugshot.bmp</p>

FIG. 2

THE UNIVERSITY OF CHICAGO

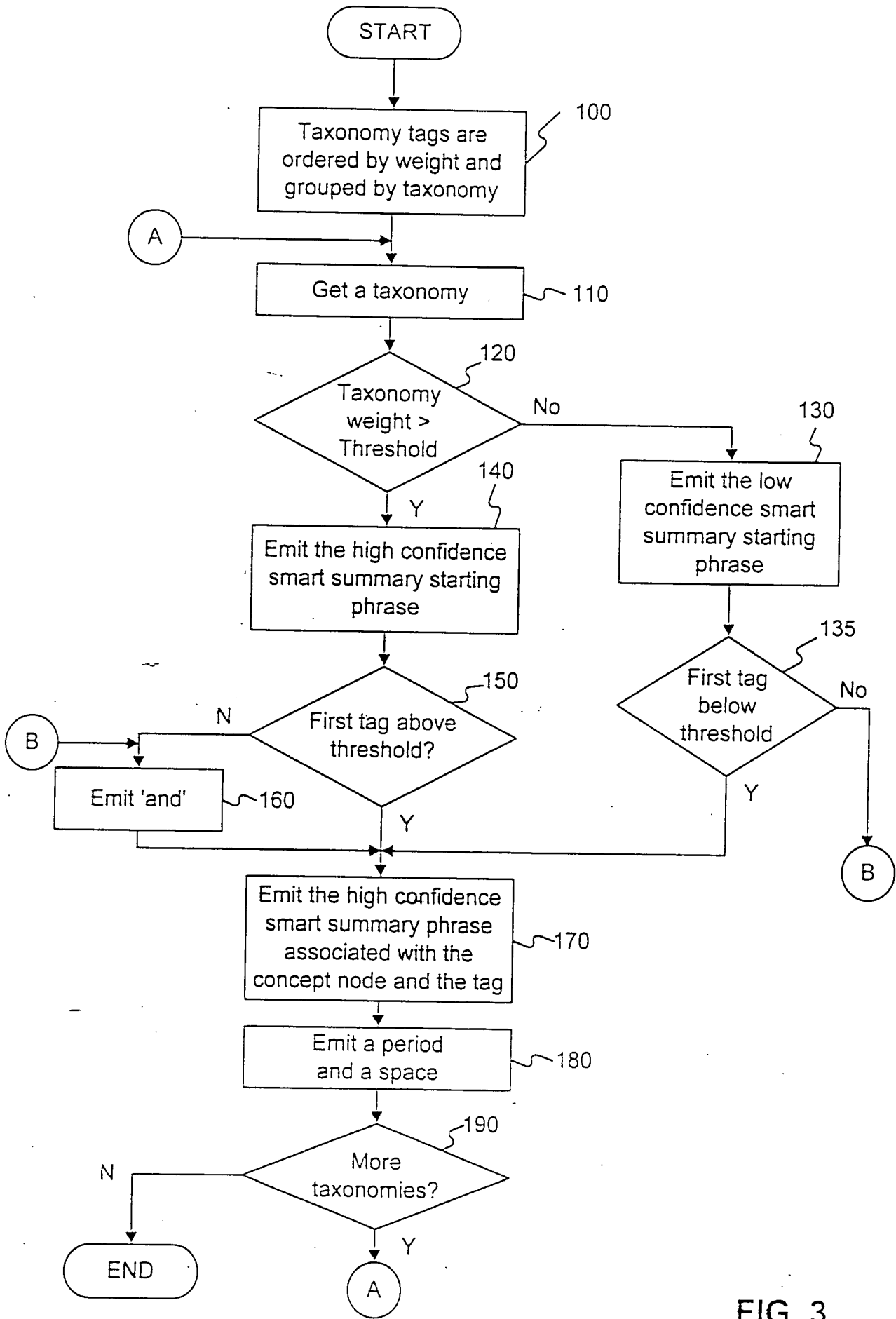


FIG. 3

Representative Taxonomy showing Types of Vehicles

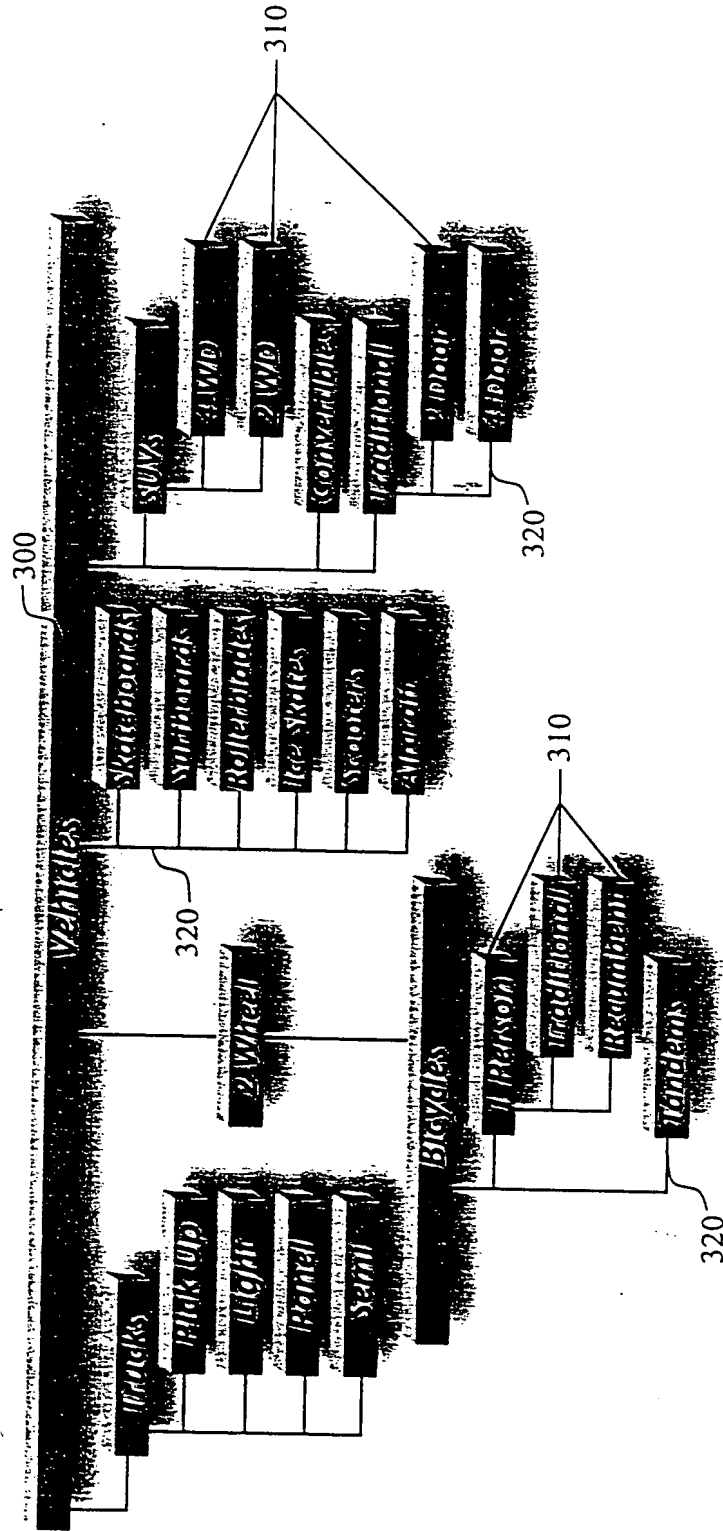


FIG. 4

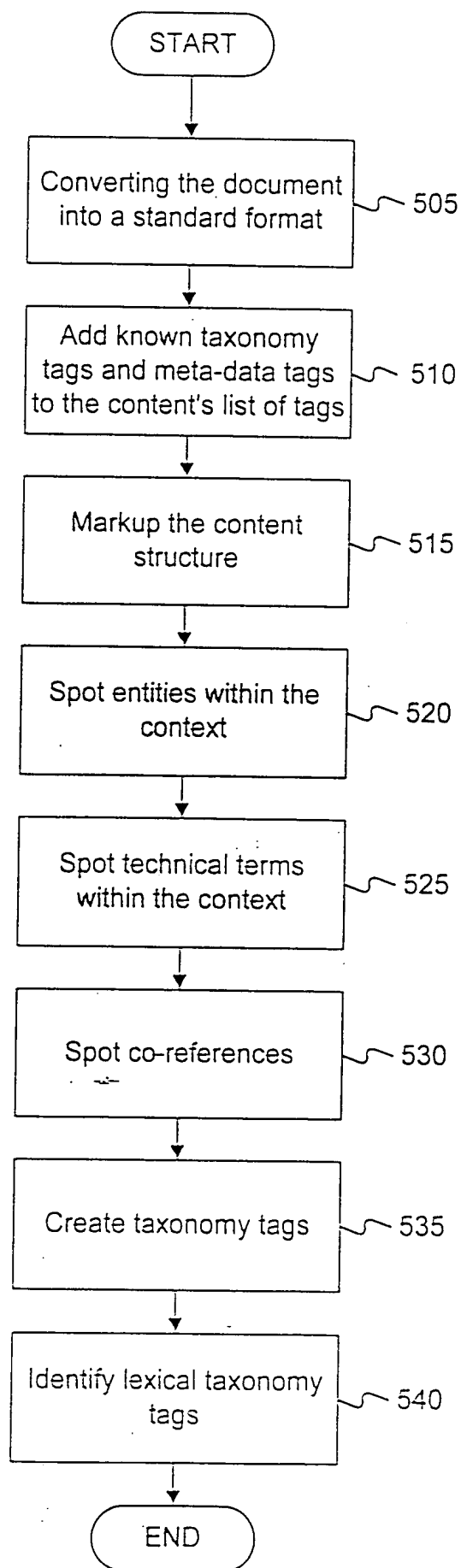
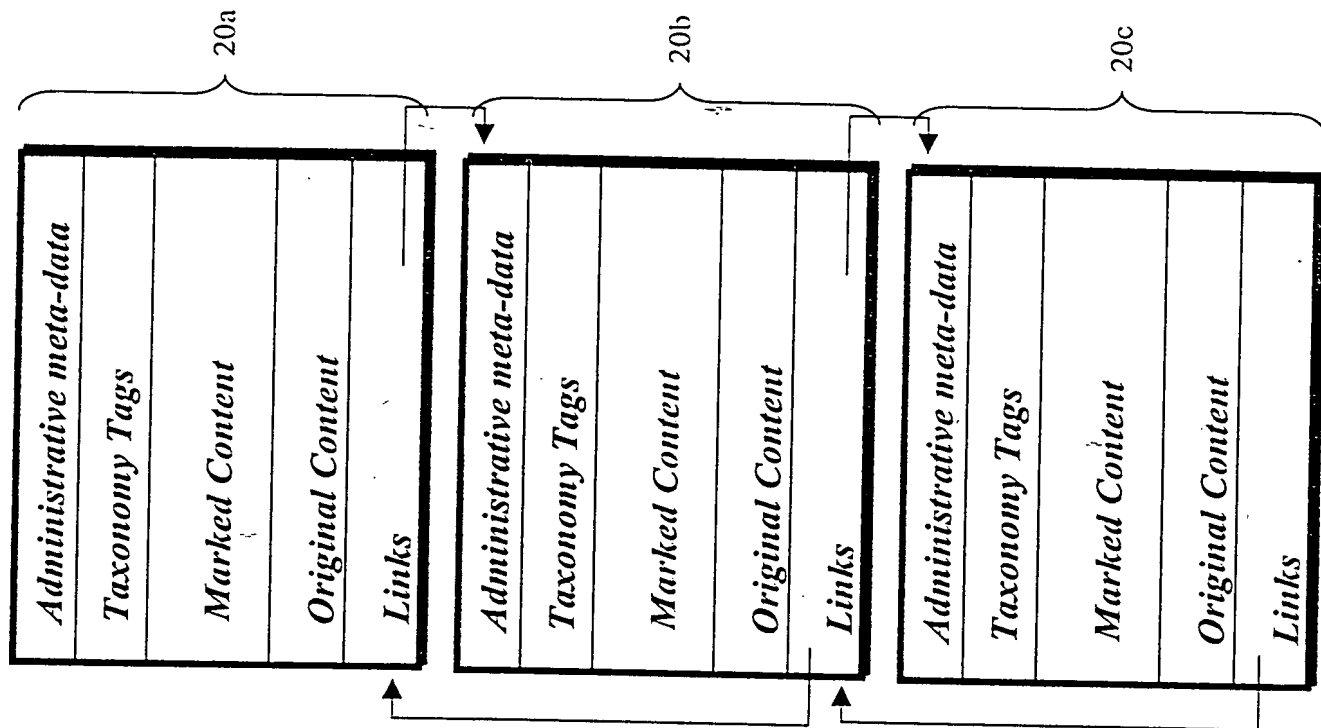
[illegible]

FIG. 5



Slice Representation

FIG. 6

FIG. 7

Slicing (2)

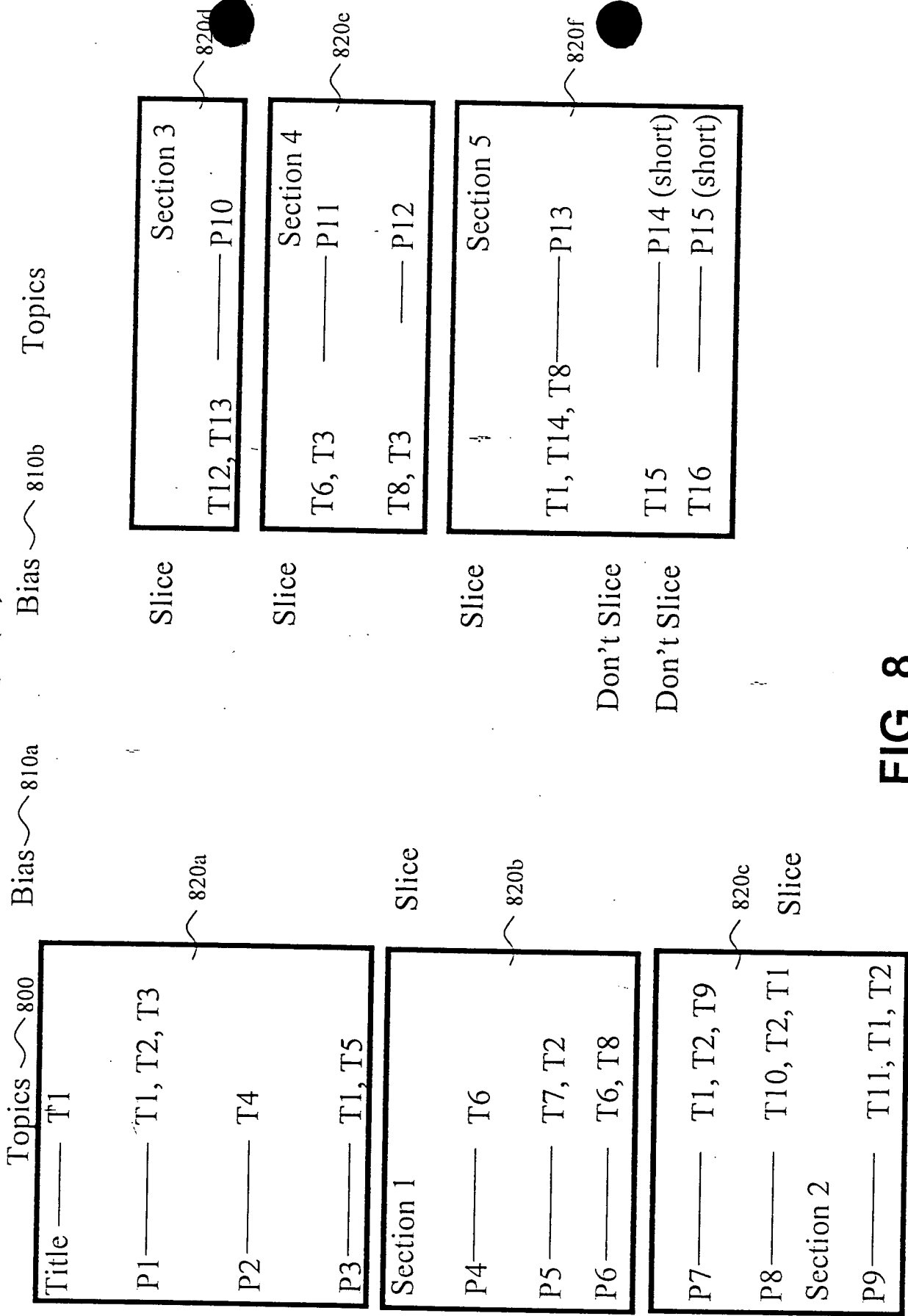
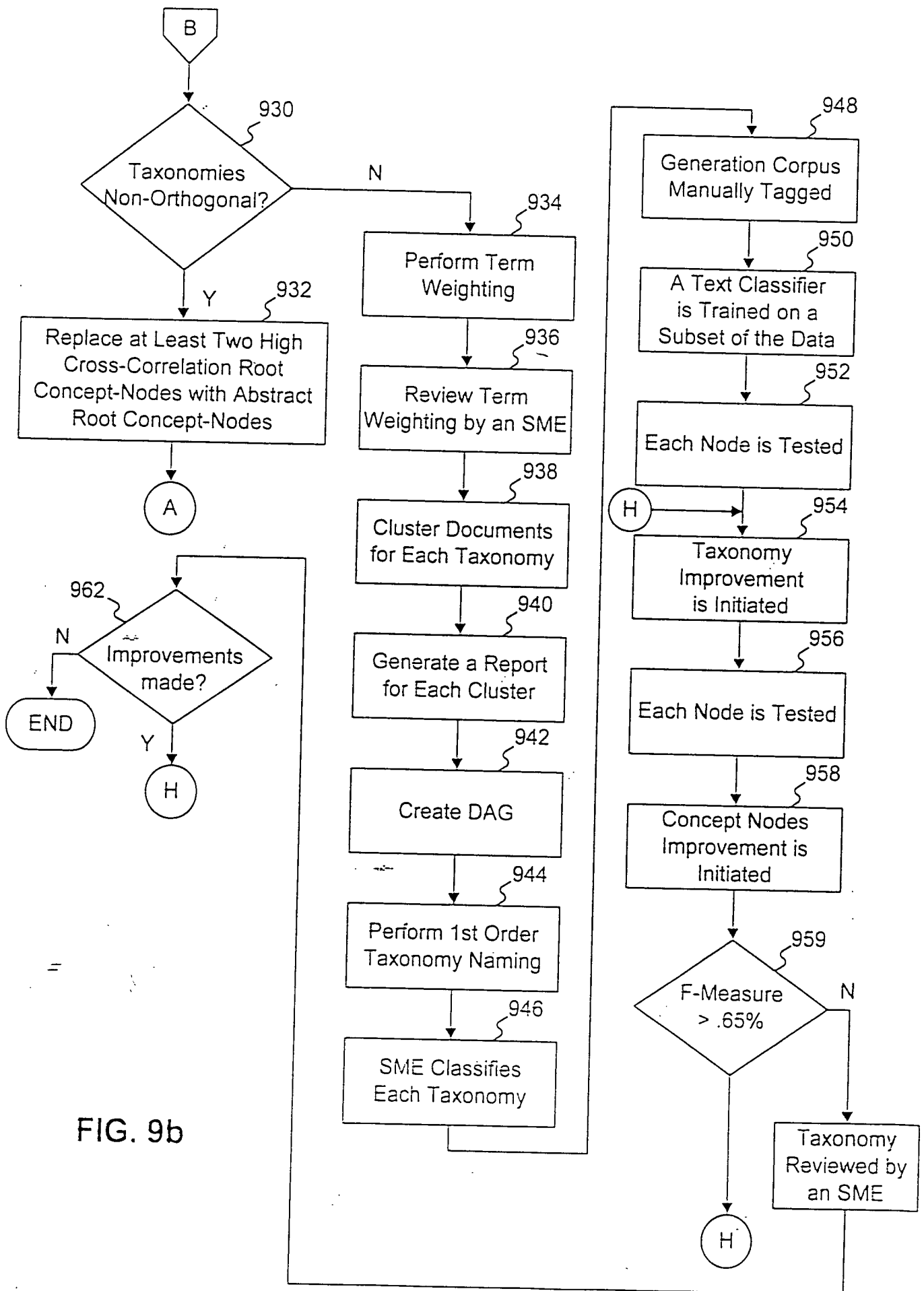


FIG. 8


```
graph TD
    START([START]) --> 902[Collect the Generation Corpus]
    902 --> 904[Collect the Taxonomy Root Concept-Nodes]
    904 --> 906[Convert the Generation Corpus into XML Documents]
    906 --> A((A))
    A --> 908[Root Concept-Node Collection and Input]
    908 --> 910[Identify and Input the Generation Corpus]
    910 --> 912[Perform Term Extraction]
    912 --> 914[Perform Term Separation]
    914 --> 916[Perform Term Analysis]
    916 --> 920[Identify Irrelevant Root Concept-Nodes]
    920 --> 922{Any Taxonomy Assigned to a Small Number of the Term/Features?}
    922 -- N --> 920
    922 -- Y --> 924[Remove Concept-Node from the Input List]
    924 --> A2((A))
    A2 --> 926{Overlap?}
    926 -- B --> B_exit{{B}}
    926 -- Y --> 928[Remove at Least One High Cross-Correlation Root Concept-Node]
    928 --> A3((A))
    A3 --> 916
```

FIG. 9a



```

graph TD
    START([START]) --> 9000[Input A, where A is a set of clusters, C]
    9000 --> 9005[Pick a cluster, C from A]
    9005 --> 9010[Identify sufficiently similar clusters, C_i]
    9010 --> 9020[Place C_i and C in partition, S]
    9020 --> 9030{More clusters in A?}
    9030 -- Y --> 9005
    9030 -- N --> 9040{S Empty?}
    9040 -- N --> 9050[Pick A cluster C in S]
    9040 -- Y --> 9045{Graph G connected with a single root?}
    9050 --> 9060[Find all clusters C_i that are similar to C]
    9060 --> 9070[Make an edge from C to every similar C_i]
    9070 --> A((A))
    9045 -- N --> D{{D}}
    9045 -- Y --> C{{C}}
  
```

FIG. 9c

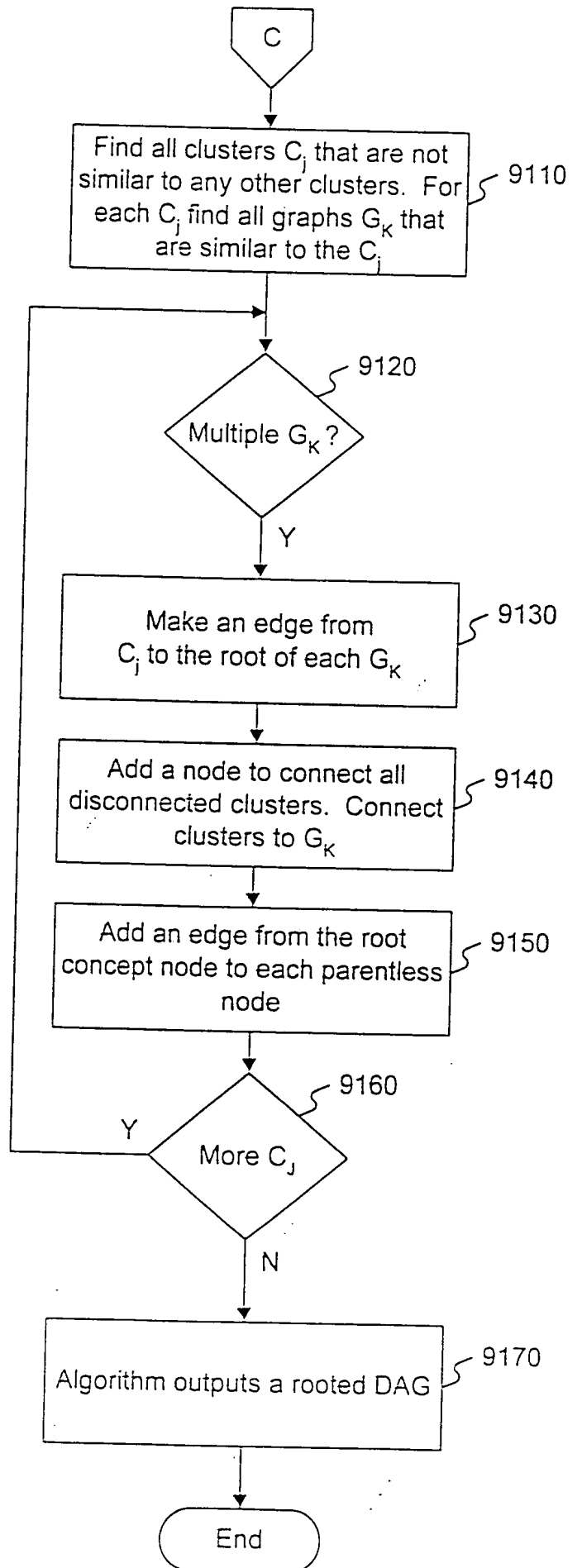
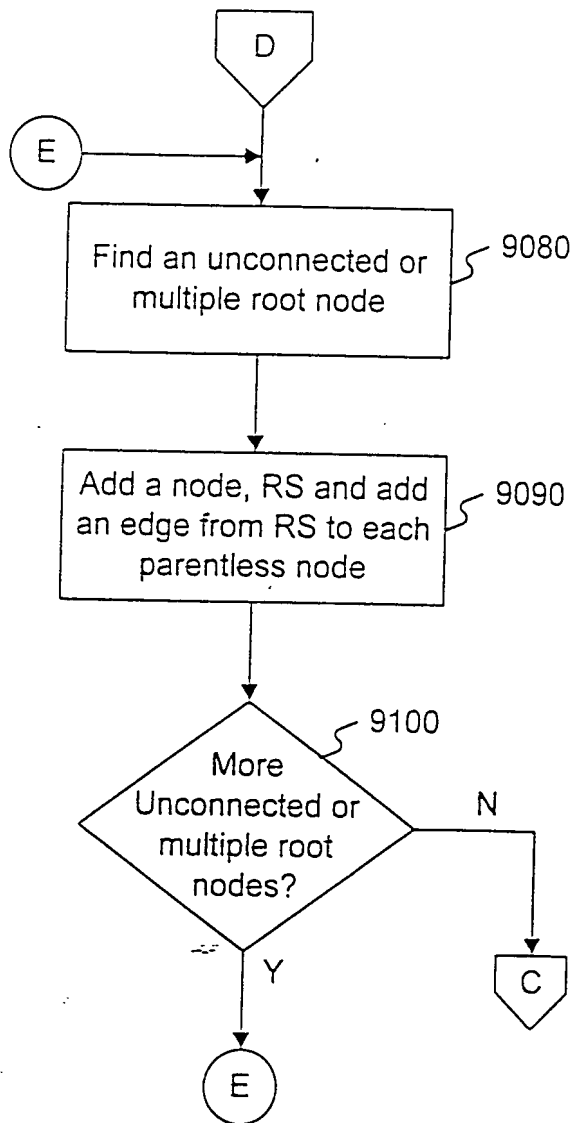


FIG. 9d

A Taxonomy of Document Sources

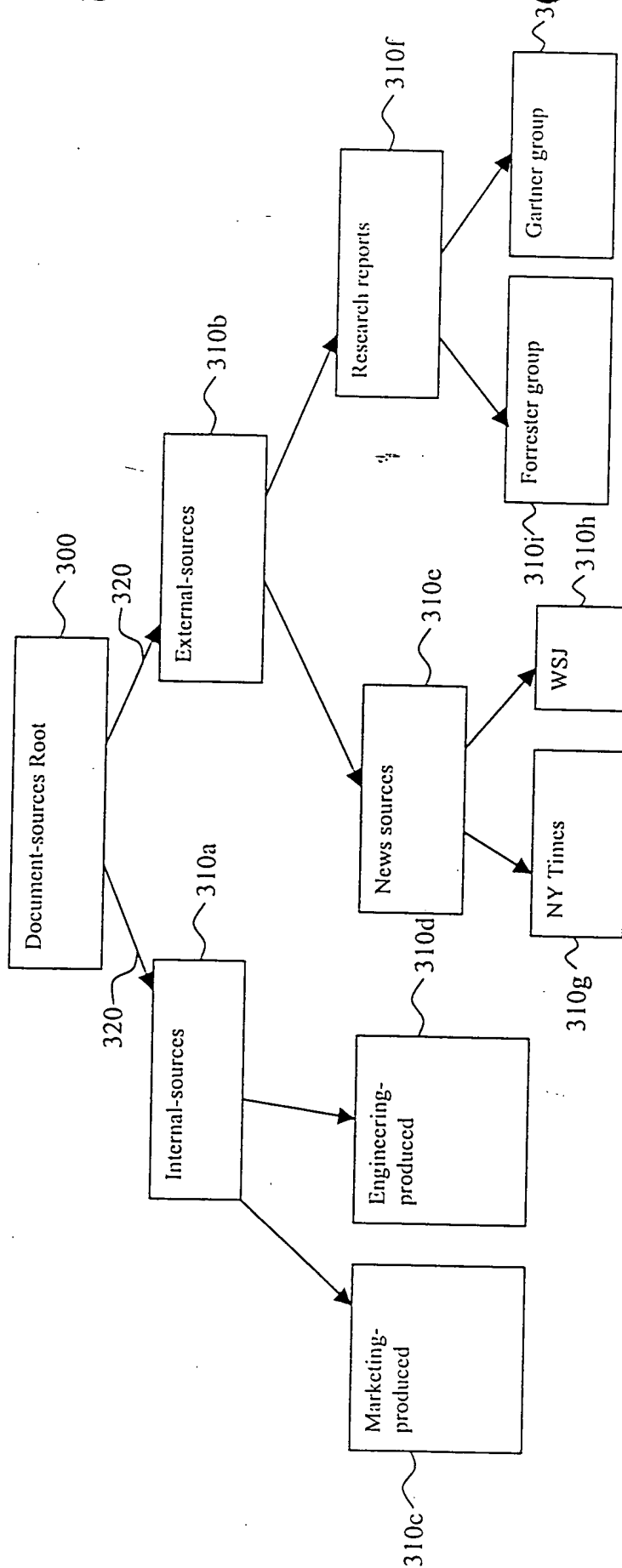
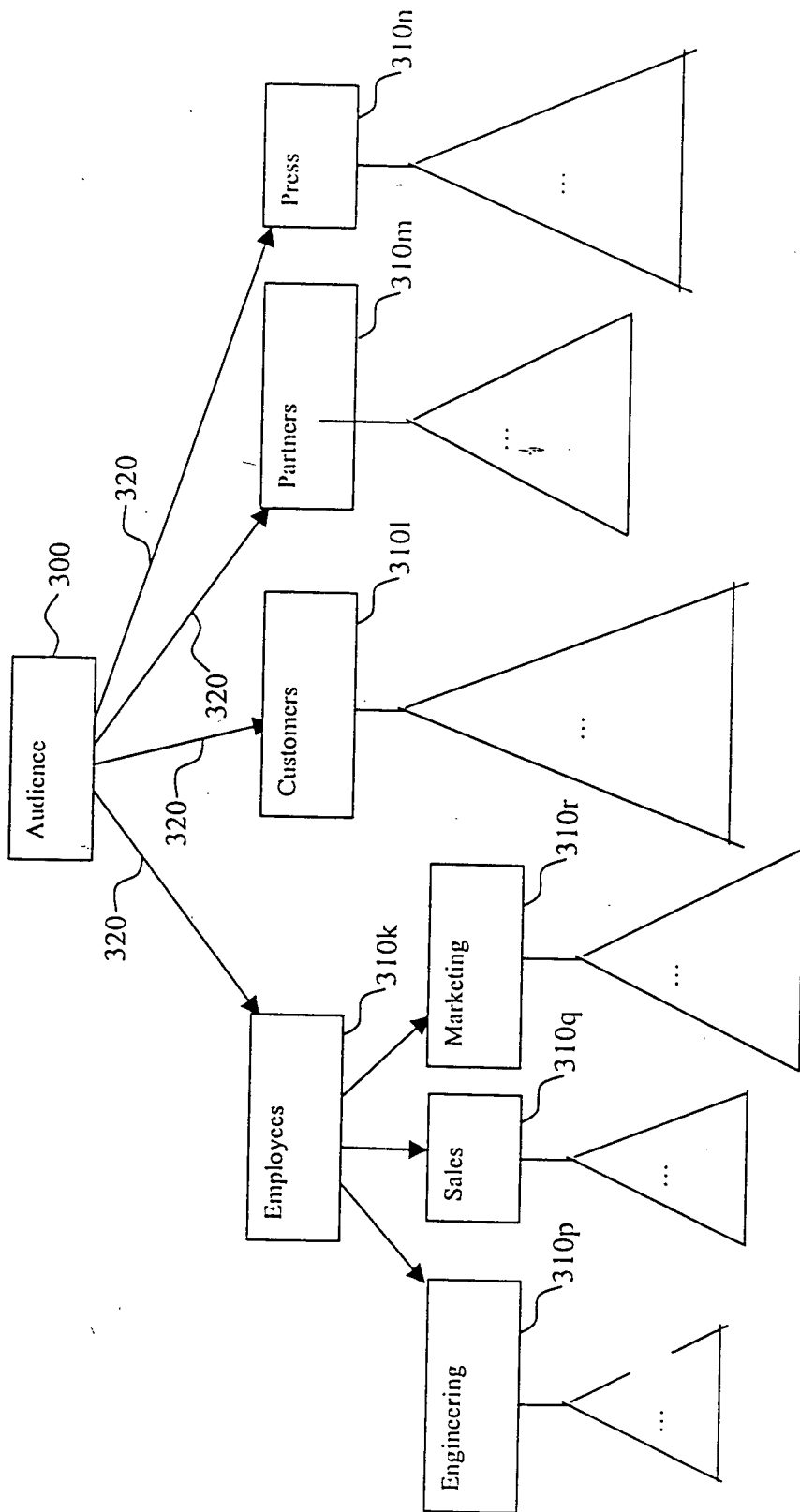


FIG. 10

Audience Taxonomy



30b

FIG. 11

Index Preparation---Step 1

For each node:

Make a "node document" which is the concatenation of all documents tagged to the node (here represented by ●).

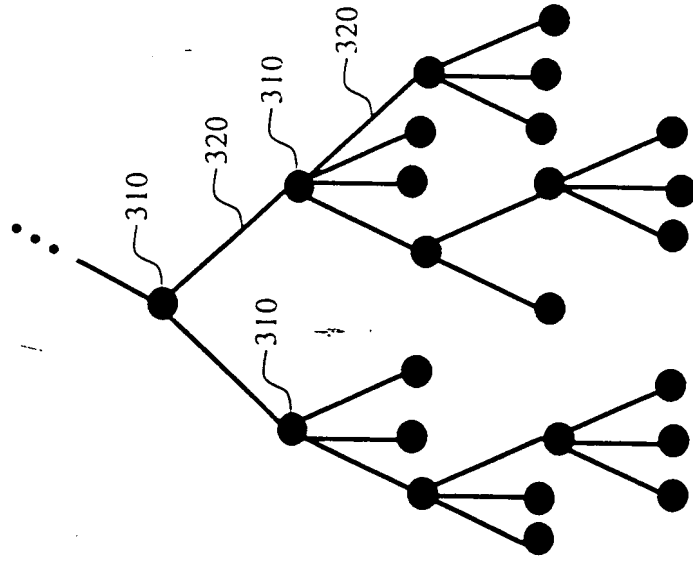






FIG. 12

Index Preparation---Step 2

Cluster the nodes according to similarity in vocabulary usage in the node-documents. Each lined-pattern represents a cluster.

-  - orange
-  - purple
-  - green
-  - blue

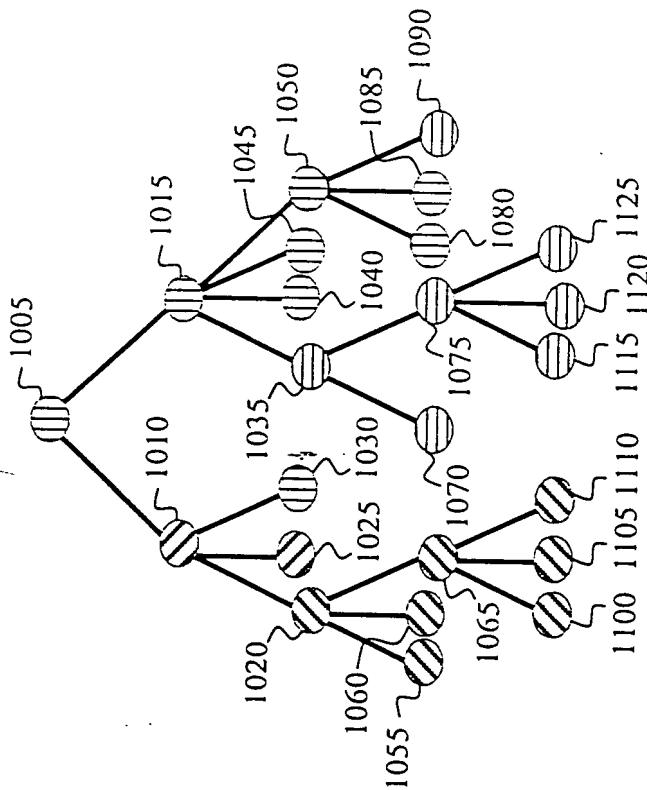
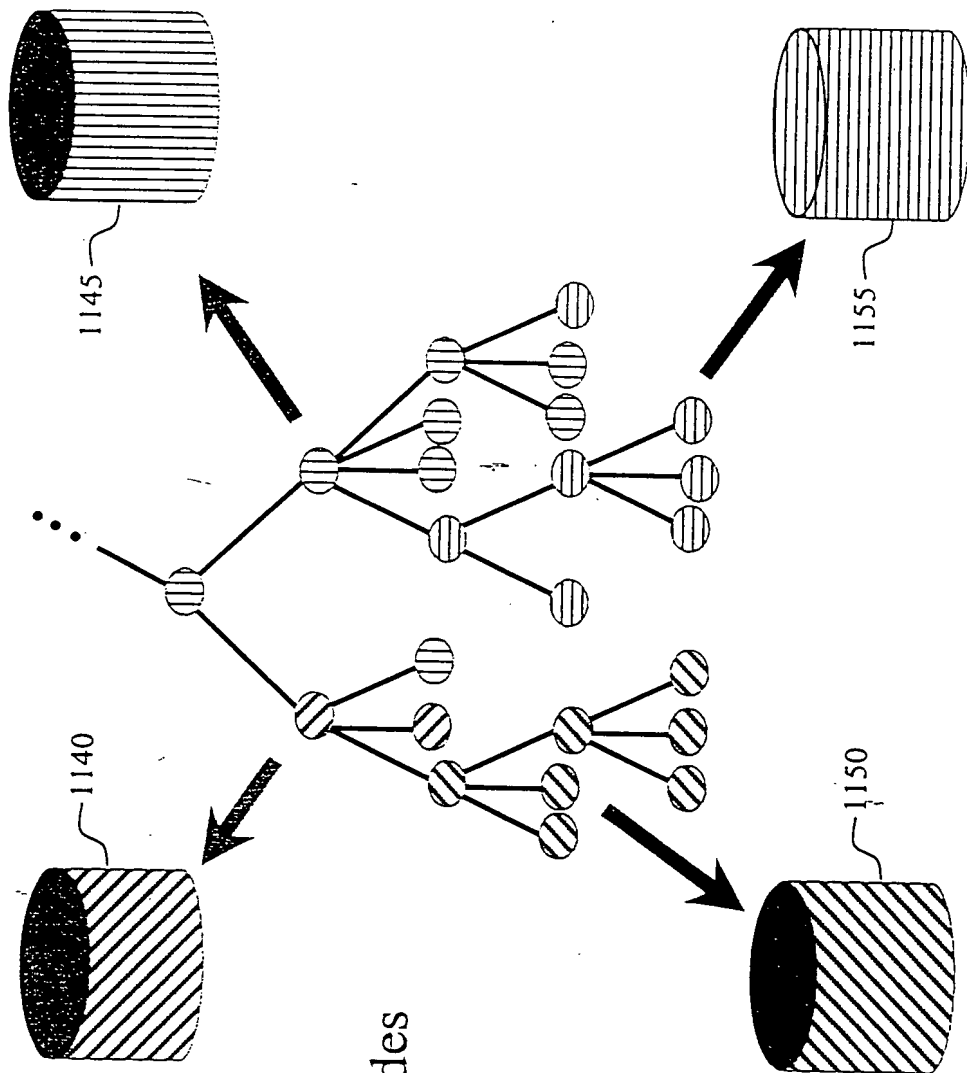


FIG. 13

Index Preparation---Step 3



For each cluster,
Build an index over the
documents tagged to nodes
in the cluster.

- orange
- purple
- green
- blue

FIG. 14

Region Designation: Marking

Identify nodes within distance D of a number N of query taxonomy tags. A node so identified is a *marked* node. Illustratively, here D is one and N is two.

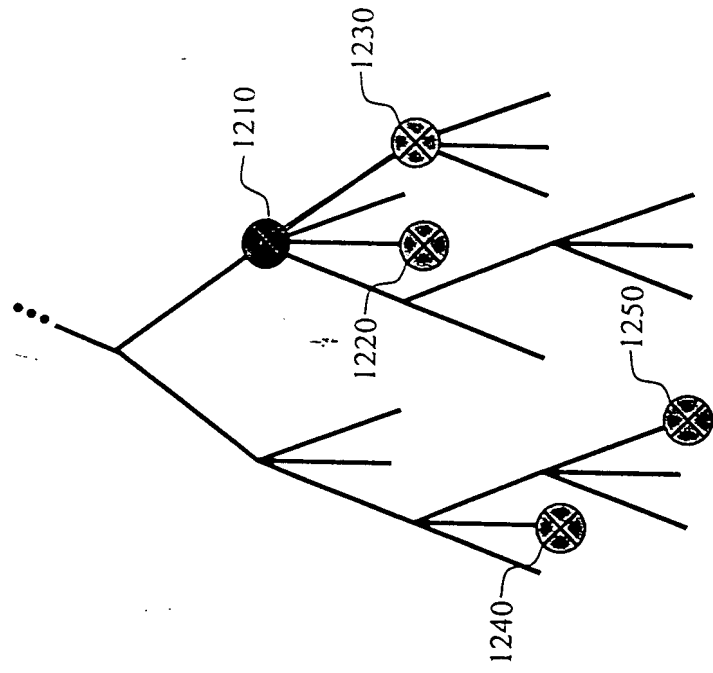


FIG. 15

Region Designation: Smoothing

Optionally, identify nodes within distance 1 of a query taxonomy tag or marked node. A node so identified is a *smoothed* node.

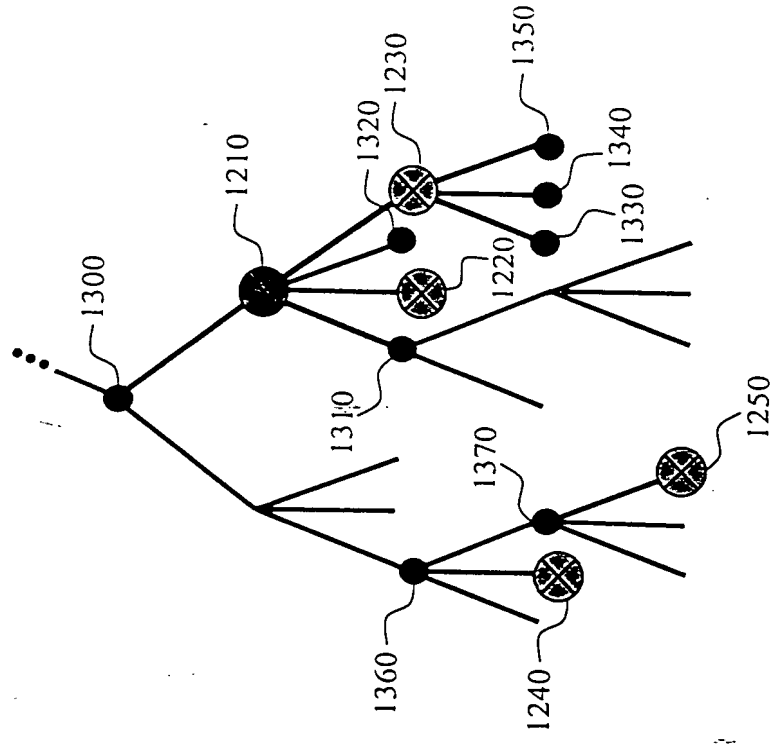


FIG. 16

Region Designation: Aggregation

Identify groups of query taxonomy tags, marked nodes, and smoothed nodes that, transitively, are within distance D' from each other. Illustratively, here D' is one.

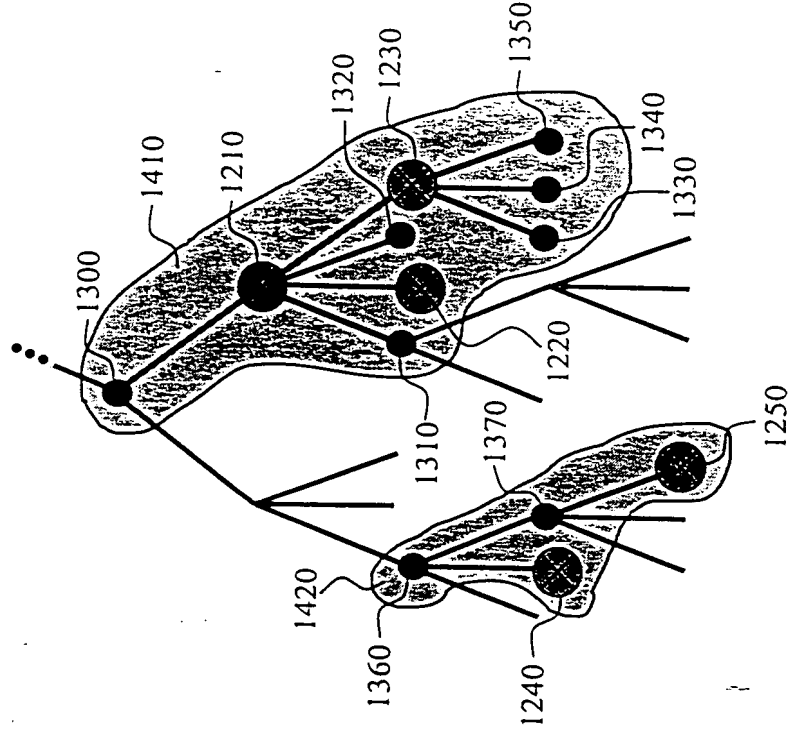






FIG. 17

Search: Index Identification

For each region,
Identify a set of indexes that
covers all of the nodes in the
region.

-  - orange
-  - purple
-  - green
-  - blue

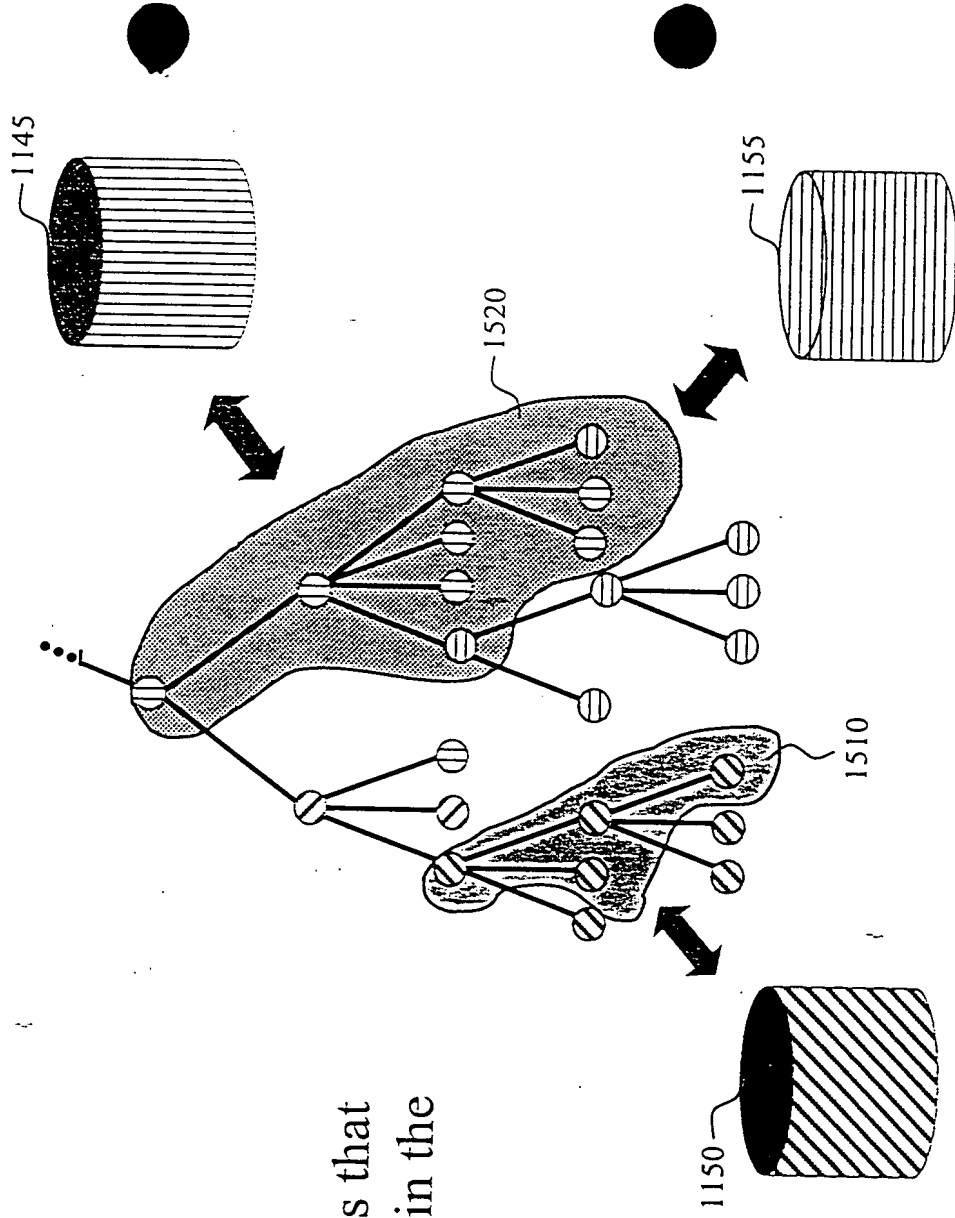


FIG. 18

Ranking: Search Engine Rank

For each document
the rank returned by
the search engine is
adjusted by ...

*For purposes of illustration,
rank is shown by distance
from the bottom of the picture.*

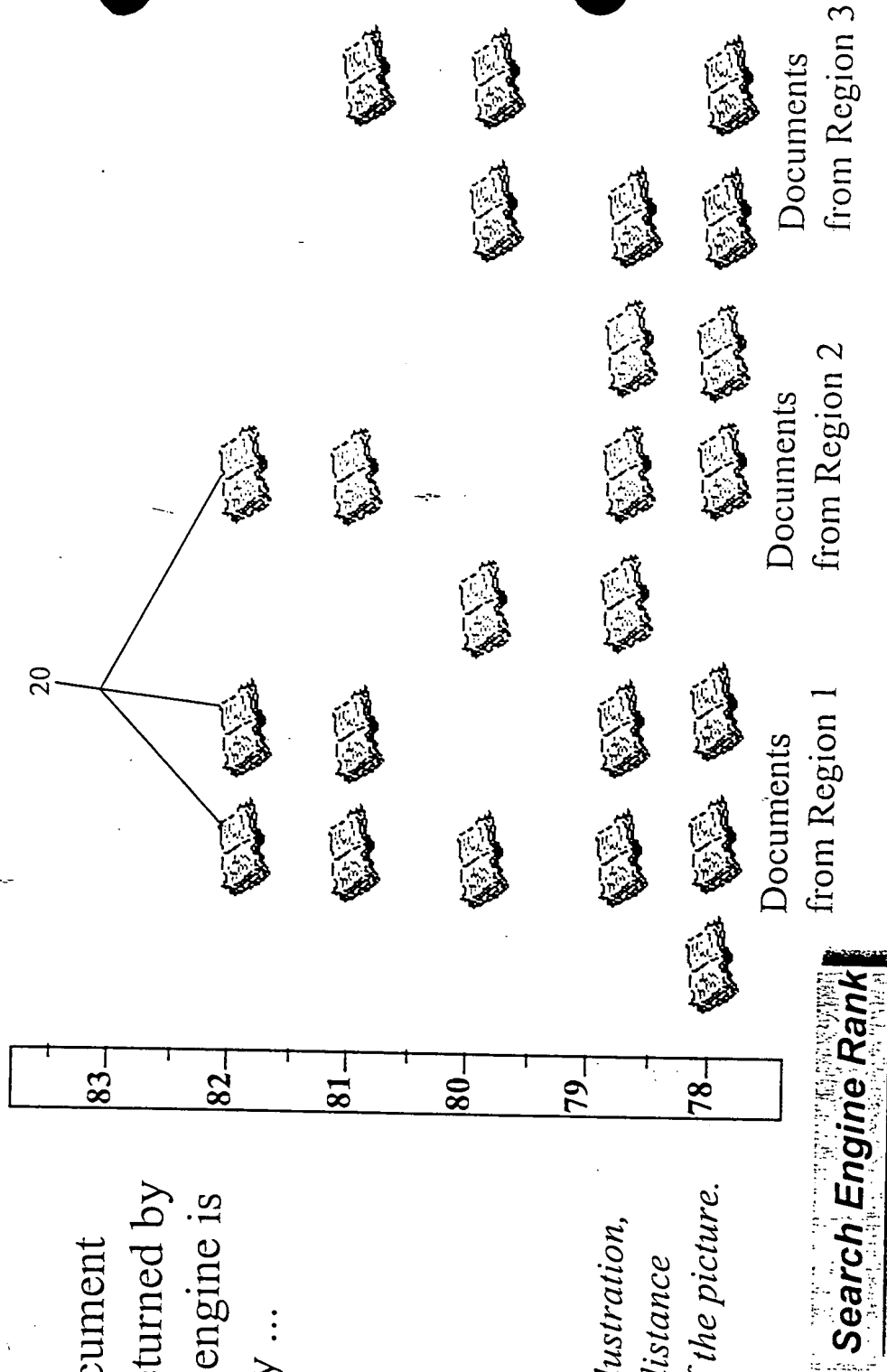
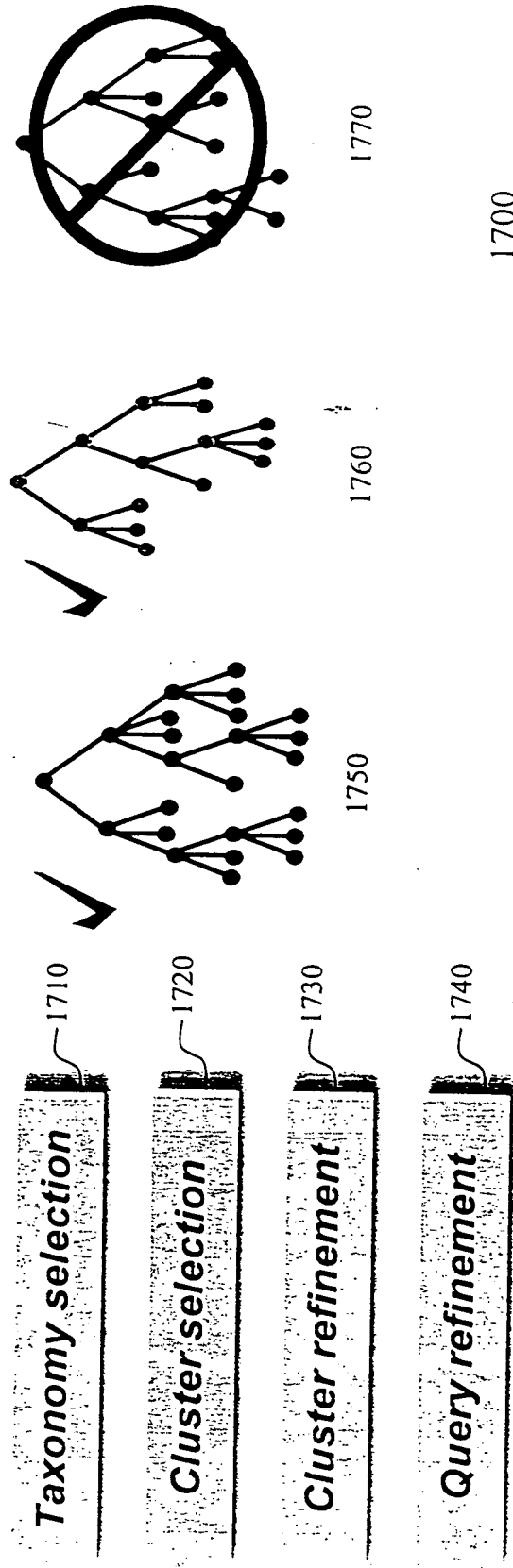


FIG. 19

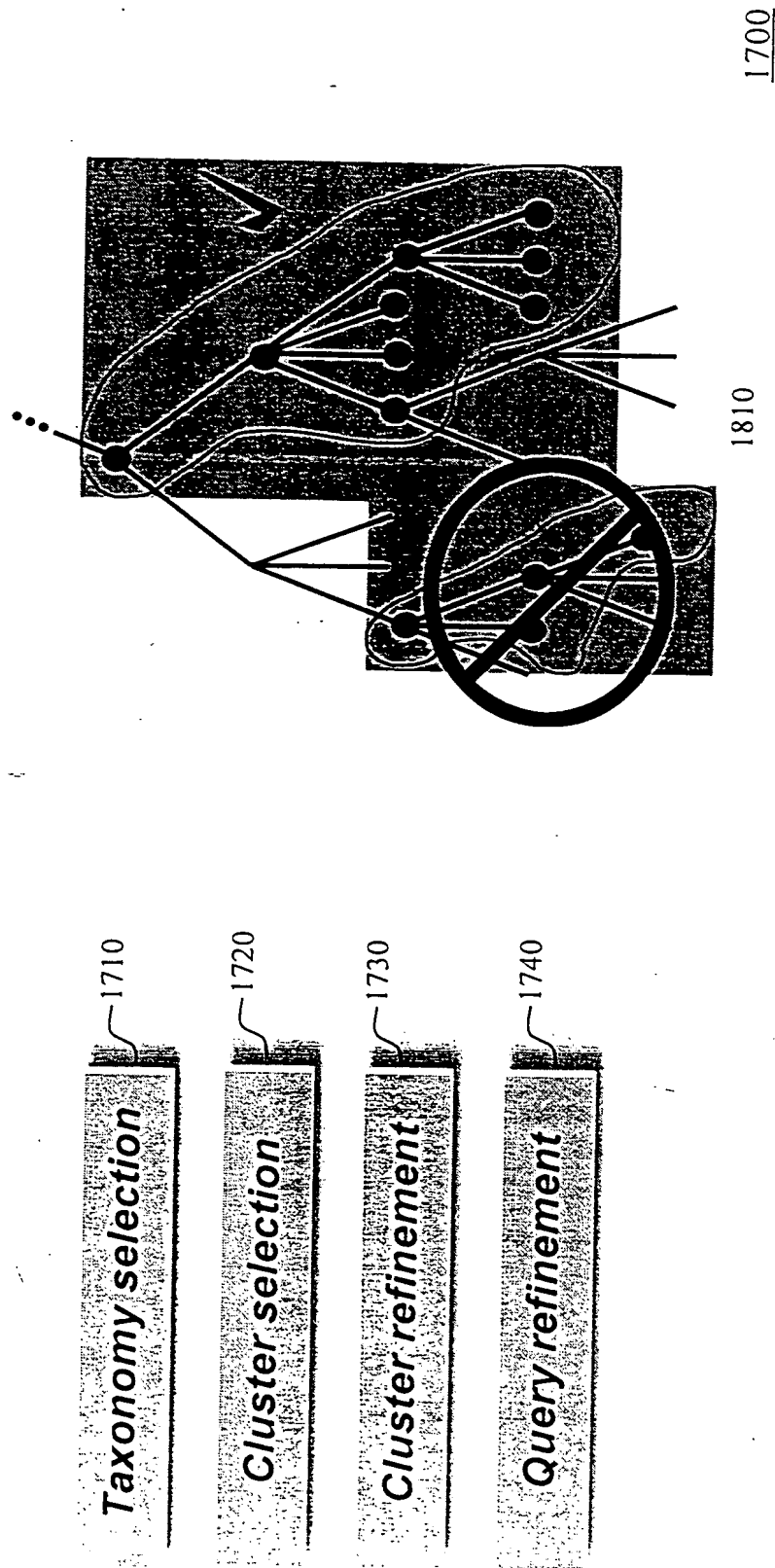
Interactive Dialogue



The user can choose among the taxonomies.

FIG. 20

Interactive Dialogue



The user can choose among the clusters.

FIG. 21

Per-topic statistics

Topic ID	Topic name	# docs	# times topic was returned	# times topic was returned correctly	Precision	Recall	F- measure
v05	Securities	52	50	49	0.98	0.94	0.96
v11	Health Care	42	42	42	1.00	1.00	1.00
v07	Manufacturing	37	38	37	0.97	1.00	0.99
v12	Public Sector	35	35	35	1.00	1.00	1.00
v01	Financial Services Overview	30	29	29	1.00	0.97	0.98
v04	Insurance	30	30	30	1.00	1.00	1.00
v10	Utilities and Energy	28	27	27	1.00	0.96	0.98
v02	Banking (Business-to-Business)	27	27	27	1.00	1.00	1.00
v03	Banking (Business-to-Consumer)	26	26	26	1.00	1.00	1.00
v08	Retail	21	23	21	0.91	1.00	0.95
v09	Service Providers	12	12	12	1.00	1.00	1.00

* Summary:

conomy has 332 documents.
 documents on average have 1.02 tags
 call = 0.99
 precision = 0.99
 measure = 0.99

Test On Train Report 1

Per-document statistics Variable style					
Document	human tags	topic spotter tags	# correct tags returned	Precision	Recall
12780	v05 (Securities)	v07 (Manufacturing)	0	0.00	0.00
12780	v01 (Financial Services Overview)	v05 (Securities)	0	0.00	0.00
12486019	v05 (Securities)	v08 (Retail)	0	0.00	0.00
18869	v10 (Utilities and Energy)	v08 (Retail)	0	0.00	0.00
18807	v05 (Securities)	v05 (Securities)	1	1.00	1.00
18807	v11 (Health Care)	v11 (Health Care)	1	1.00	1.00
12619	v03 (Banking (Business-to-Consumer))	v03 (Banking (Business-to-Consumer))	1	1.00	1.00
18807	v01 (Financial Services)	v01 (Financial Services)	1	1.00	1.00

Test On Train Report 2

(v1.0)	Discriminating terms for v05 - Securities
2.42	Charles Schwab Securities Industry
1.77	Tax Relief Act
1.65	Securities and Exchange Commission
1.46	Merrill Lynch
1.42	Charles Schwab Corp.
1.41	compound annual growth rate
1.40	full-service firm
1.37	brokerage firm
1.37	minimum account balance
1.37	industry-wide testing
1.37	industrywide testing
1.33	discount brokerage
1.32	online brokerage
1.30	Ameritrade Holding Corporation
1.27	Wall Street
1.26	online broker
1.25	external resource
1.21	securities firm
1.19	IT staff
1.16	capacity expansion
1.16	technology development
1.10	implementation concern
1.09	Capacity Study
1.09	Securities Industry Association
1.08	mutual fund
1.08	infrastructure project
1.02	Tower Group
1.01	National Securities Clearing Corporation
1.00	staffing resource

[Edit](#) [View](#) [Go](#) [Fav](#) [Printer](#) [Help](#)
[Back](#) [Forward](#) [Stop](#) [Refresh](#) [Home](#) [Search](#) [Fav](#) [History](#) [Unblock](#) [Full Screen](#) [Mail](#) [Print](#) [Edit](#)
 (98%) C:\Cisco\Taxonomies\allisontaxos\CiscoVerticals\VariableTestOnTrain\doc_lists\v05_ht.html

Topic: v05 (Securities)
 op terms used by topic spotter:

1. Charles Schwab Securities Industry
2. Tax Relief Act
3. Securities and Exchange Commission
4. Merrill Lynch
5. Charles Schwab Corp.

Document	Classification	Human-assigned tags	Topic spotter tags	Topic spotter scores
18667		v05	v05 (Securities)	0.25
18668		v05	v05 (Securities)	0.35
18669		v05	v05 (Securities)	0.34
18670		v05	v05 (Securities)	0.20
18671		v05	v05 (Securities)	0.27
18672		v05	v05 (Securities)	0.18

Test On Train Report 4

C:\Cisco\Taxonomies\allisontaxos\Cisco\Verticals\VariableTestOnTrain\doc_lists\v05_hl.html				
128security		v05	v05 (Securities)	0.18
124security	should be tagged to v05	v05	v08 (Retail)	0.13
12755	should be tagged to v05	v05	v07 (Manufacturing)	0.25
1255security		v05	v05 (Securities)	0.28
001\03\05\02\security\0999	should be tagged to v05	v01 v03 v05	v03 (Banking (Business-to-Consumer)) v01 (Financial Services Overview)	0.25 0.23
1268security		v05	v05 (Securities)	0.21
127security		v05	v05 (Securities)	0.21
1288security		v05	v05 (Securities)	0.24
1298security		v05	v05 (Securities)	0.30
1300		v05	v05 (Securities)	0.32
1301		v05	v05 (Securities)	0.19

Test On Train Report 5